

Systèmes du LIA à DEFT'13

Xavier Bost¹ Ilaria Brunetti^{1,6} Luis Adrián Cabrera-Diego¹
 Jean-Valère Cossu¹ Andréa Linhares^{1,5} Mohamed Morchid¹
 Juan-Manuel Torres-Moreno^{1,2,3,4} Marc El-Bèze^{1,2,4} Richard
 Dufour^{1,4}

(1) LIA, 339, chemin des Meinajariès 84912 Avignon Cedex 09

(2) SFR Agorantic Université d'Avignon et des Pays de Vaucluse, 84000 Avignon Cedex

(3) École Polytechnique de Montréal, 2900 Bd Edouard-Montpetit Montréal, QC H3T1J4

(4) Brain & Language Research Institute, 5 avenue Pasteur, 13604 Aix-en-Provence Cedex 1

(5) Universidade Federal do Ceara Rua Estanislau Frota, S/N, CEP 62.010-560 – Sobral/ CE
 Brésil

(6) Equipe Maestro INRIA, 2004, Route des Lucioles, 06902 Sophia Antipolis
 prenom.nom@univ-avignon.fr

RÉSUMÉ

La campagne Défi de Fouille de Textes (DEFT) en 2013 s'est intéressée à deux types de fonctions d'analyse du langage, la classification de documents et l'extraction d'information dans le domaine de spécialité des recettes de cuisine. Nous présentons les systèmes du LIA appliqués à DEFT 2013. Malgré la difficulté des tâches proposées, des résultats intéressants ont été obtenus par nos systèmes.

ABSTRACT

The Systems of LIA at DEFT'13

The 2013 *Défi de Fouille de Textes* (DEFT) campaign is interested in two types of language analysis tasks, the document classification and the information extraction in the specialized domain of cuisine recipes. We present the systems that the LIA has used in DEFT 2013. Our systems show interesting results, even though the complexity of the proposed tasks.

MOTS-CLÉS : Classification de textes ; Extraction d'information ; Méthodes statistiques ; Domaine de spécialité.

KEYWORDS: Text Classification ; Information Extraction ; Statistical Methods ; Speciality domain.

1 Description des tâches et des données

Le Défi Fouille de Textes (DEFT) s'intéresse en 2013 à deux types de fonction d'analyse du langage, la classification de documents (tâches 1 à 3) et l'extraction d'information (tâche 4), ceci dans un domaine de spécialité : les recettes de cuisine. Pour cette nouvelle édition du défi¹, l'équipe du Laboratoire Informatique d'Avi-

1. <http://deft.limsi.fr/2013/>

gnon (LIA) a décidé de participer aux tâches suivantes :

- T1. L'identification à partir du titre et du texte de la recette de son niveau de difficulté sur une échelle à 4 niveaux : *très facile, facile, moyennement difficile, difficile* ;
- T2. L'identification à partir du titre et du texte de la recette du type de plat préparé : *entrée, plat principal, dessert* ;
- T4. L'extraction à partir du titre et du texte d'une recette de la liste de ses ingrédients.

Le corpus, fourni par DEFT au format XML, est composé de 13 684 recettes. Nous avons choisi d'augmenter notre base de données en extrayant des recettes de cuisine à partir de 4 sites spécialisés : Cuisine AZ (<http://www.cuisineaz.fr>), Madame Le Figaro (<http://madame.lefigaro.fr/recettes>), 750 grammes (<http://www.750g.com/>) et Ptit Chef (<http://www.ptitchef.com/>). Environ 50k recettes supplémentaires ont pu être obtenues à partir de ces sites. Notons cependant que les recettes téléchargées ne comportaient pas toujours les mêmes champs d'information que ceux contenus dans les recettes fournies par les organisateurs de la campagne.

2 Méthodes utilisées

Pour répondre à ce défi, nous avons employé plusieurs méthodes de classification (modèles probabilistes discriminants, *boosting* de caractéristiques textuelles) et d'extraction d'information, ainsi que leurs combinaisons.

2.1 Modèles discriminants

Le système DISCOMP_LIA a été appliqué aux tâches 1 et 2. Son fonctionnement repose essentiellement sur les modèles discriminants décrits dans (Torres-Moreno et al., 2011). Ces modèles servent aussi bien à agglutiner des mots afin de produire des termes composites ayant un fort pouvoir discriminant (sans toutefois passer en dessous d'un seuil préfixé de couverture), qu'à pondérer de façon plus appropriée ces termes dans le calcul de l'indice de similarité (ici Cosine revisité) entre chaque recette et le sac de mots de chacune des classes (3 ou 4 selon les tâches).

Trois nouveautés ont été introduites pour prendre en compte les particularités du défi :

- nous avons opté pour une approche hiérarchique propre à chacune des 2 tâches :
 - pour la tâche T1 : identification du niveau *facile* ou *difficile* d'une recette, puis dans chacun des 2 groupes, différenciation en 2 sous-groupes (*moyennement* ou *très difficile*) ;
 - pour la tâche T2 : identification du type *dessert* ou *autre* ; dans le cas d'une identification en type *autre*, différenciation entre *plat principal* et *entrée* ;
- nous avons appliqué le même système avec un paramétrage *ad hoc* à 2 jeux

de données : le titre de la recette seul pour le premier et pour le second l'ensemble formé par le contenu et le titre de la recette. Une combinaison linéaire entre les scores des 2 systèmes a permis d'obtenir une nouvelle distribution de probabilités sur les hypothèses de classes ;

- la formation par agglutination des termes et ce, surtout dans le titre, a fait apparaître des sous-types de recettes (quiches, soupes, gâteaux) qui ont permis fournir aux modèles de meilleurs points d'appui pour identifier la classe visée. Dans les cas où l'indice de pureté de *Gini* valait 1, et où la couverture était élevée, nous avons enrichi le sac de mots avec des composants factices de façon à compenser le bruit qui pouvait être engendré par le reste de la recette. Ceci a été fait de façon semi-automatique, mais devrait pouvoir être complètement automatisé.

Notons par ailleurs que nous nous sommes servis du critère de pureté de *Gini* pour sélectionner quelques mots erronés ce qui, après correction, a permis d'augmenter leur couverture sans altérer leur pouvoir discriminant.

2.2 Algorithmes de *boosting*

Pour les tâches de classification (1 et 2), nous proposons une manière de combiner différentes caractéristiques textuelles et numériques au moyen d'un algorithme d'apprentissage automatique : le *Boosting* (Schapire, 2003). Son objectif principal est d'améliorer (de « booster ») la précision de n'importe quel algorithme d'apprentissage permettant d'associer, à une série d'exemples, leur classe correspondante. Le principe général du *Boosting* est assez simple : la combinaison pondérée d'un ensemble de classifieurs binaires (appelés classifieurs *faibles*), chacun associé à une règle différente de classification très simple et peu efficace, permettant au final d'obtenir une classification robuste et très précise (classifieur *fort*). La seule contrainte de ces classifieurs *faibles* est d'obtenir des performances meilleures que le hasard. Dans le cas des classifieurs binaires, l'objectif est de classer plus de 50 % des données correctement.

Nous proposons d'utiliser l'algorithme *AdaBoost* car il présente de nombreux avantages dans le cadre des tâches de classification proposées. En effet, l'algorithme est très rapide et simple à programmer, applicable à de nombreux domaines, adaptable aux problèmes multi-classes. Il fonctionne sur le principe du *Boosting*, où un classifieur *faible* est obtenu à chaque itération de l'algorithme *AdaBoost*. Chaque exemple d'apprentissage possède un poids, chaque tour de classification permettant de les re-pondérer selon le classifieur *faible* utilisé. Le poids d'un exemple bien catégorisé est diminué, au contraire d'un exemple mal catégorisé, pour lequel on augmente son poids. L'itération suivante se focalisera ainsi sur les exemples les plus « difficiles » (poids les plus élevés). L'algorithme a été implémenté dans l'outil de classification à large-marge *BoosTexter* (Schapire and Singer, 2000), permettant de fournir comme caractéristiques au classifieur des données numériques mais également des données textuelles. Les classifieurs *faibles* sur les données textuelles peuvent prendre en compte des *n*-grammes de mots. Nous avons utilisé l'outil *IcsiBoost* (Favre et al., 2007), qui est l'implémentation open-source de cet outil. *IcsiBoost* a notamment l'avantage de pouvoir fournir, pour chaque exemple à classer, un score de confiance pour chacune des classes entre 0 (peu confiant) et 1 (très confiant).

2.2.1 Niveau de difficulté

Pour l'approche utilisant l'algorithme de classification *AdaBoost*, nous avons retenu différentes caractéristiques qui ont été fournies en entrée pour l'apprentissage du classifieur :

- données textuelles
 1. n -grammes de mots du titre de la recette avec taille maximum de 3 mots ;
 2. n -grammes de mots du texte de la recette avec taille maximum de 4 mots ;
 3. uni-gramme sur la liste des ingrédients obtenue automatiquement (voir section 2.6).
- données numériques continues
 1. nombre de mots du titre de la recette ;
 2. nombre de mots du texte de la recette ;
 3. nombre de phrases du texte de recette ;
 4. nombre de séparateurs du texte de la recette (point, virgule, deux-points, etc.) ;
 5. taille de la liste des ingrédients obtenue automatiquement.

À la fin de ce processus d'apprentissage, la liste des classifieurs sélectionnés est obtenue tout comme le poids de chacun d'eux, afin d'utiliser les exemples les plus discriminants pour chaque classe (chaque niveau de difficulté est considéré dans cette tâche comme une classe). Nous avons composé un sous-corpus à partir du corpus d'apprentissage fourni par les organisateurs de la tâche afin d'optimiser le nombre de classifieurs *faibles* nécessaires. Ce sous-corpus est composé de 3 863 recettes respectant la distribution des classes, le reste étant utilisé pour l'apprentissage (environ 72% des données d'apprentissage). Notons que tout le corpus d'apprentissage sera néanmoins utilisé pour entraîner les classifieurs pendant la phase de test de la campagne. Enfin, une attention particulière a été portée sur la normalisation des données textuelles du titre et de la description des recettes. En effet, les données textuelles « brutes » fournies par les organisateurs possédaient un vocabulaire très bruité principalement à cause des nombreuses abréviations (par exemple « th » pour « thermostat ») et fautes d'orthographe (par exemple « échalote » et « échalotte »). Quatre traitements particuliers ont été appliqués sur le texte :

1. suppression de la ponctuation et isolation de chaque mot (exemple : « et couper l'oignon. » devient « etLcouperLl'l'oignon ») ;
2. utilisation d'une liste manuelle de normalisation de mots (exemple : « kg » devient « kilogramme ») ;
3. conversion automatique des chiffres en lettres ;
4. regroupement des n -grammes de mots les plus fréquents au sein d'une même entité au moyen de l'outil *lia_tagg*² (exemple : « il y a » est considéré comme un mot « il_y_a »).

Ce texte normalisé est utilisé par l'algorithme *AdaBoost* pour la recherche des classifieurs *faibles* sur les données textuelles lors de la phase d'apprentissage, mais également pour la recherche du niveau de difficulté associé à une recette lors de la

2. http://lia.univ-avignon.fr/fileadmin/documents/Users/Intranet/chercheurs/bechet/download_fred.html

phase de test. Au final, nous obtenons pour chaque recette à classer, un score de confiance pour chaque niveau de difficulté, permettant de choisir le niveau de difficulté de la recette (score le plus élevé).

2.2.2 Type de plat

Dans ce problème de classification au moyen de l'approche par *Boosting*, nous avons utilisé les mêmes caractéristiques et normalisations que celles décrites dans la tâche de recherche du niveau de difficulté (voir section 2.2.1). Pour chaque recette, nous obtenons donc un score de confiance pour chaque type de plat qui pourra être utilisé pour choisir la classe à associer à une recette.

2.3 SVMs

Le corpus des recettes \mathbb{X} étant donné, il s'agit d'élaborer à partir d'un sous-ensemble $\mathbb{D} \subset \mathbb{X}$ de recettes annotées par classe (niveau de difficulté ou type de plat selon la tâche), une méthode de classification γ qui à toute recette associe une classe. En notant \mathbb{C} l'ensemble des classes, nous avons donc :

$$\gamma : \mathbb{X} \longrightarrow \mathbb{C} \quad (1)$$

La troisième des méthodes appliquées a consisté, pour chaque couple de classes $(c, c') \in \mathbb{C} \times \mathbb{C}$ (avec $c \neq c'$), à apprendre et à appliquer un classifieur binaire $\gamma_{(c, c')} : \mathbb{X} \rightarrow \{c, c'\}$ à base de SVMs.

En notant $s_{(c, c')}(r)$ le score obtenu selon le classifieur $\gamma_{(c, c')}$ par la recette $r \in \mathbb{X}$, nous avons alors, avec $c' \in \mathbb{C}$:

$$\gamma(r) = \arg \max_{c \in \mathbb{C}} (s_{(c, c')}(r)) \quad (2)$$

Lors des phases d'apprentissage et d'application des classifieurs binaires, les recettes ont été représentées vectoriellement dans l'espace du lexique de référence. Le poids $w(t)$ du terme t d'une recette r dans le vecteur associé est donné par :

$$w(t) = tf(t).idf(t) = tf(t) \cdot \log \left(\frac{|\mathbb{X}|}{df(t)} \right) \quad (3)$$

où $tf(t)$ désigne le nombre d'occurrences du terme t dans la recette et $df(t)$ le nombre de recettes du corpus \mathbb{X} où le terme t apparaît. L'approche décrite a été appliquée pour les deux tâches de classification.

Dans le cas de la prédiction du type de plat, le lexique a cependant été filtré avant vectorisation des recettes sur le critère de l'information mutuelle entre un terme t et une classe $c \in \mathbb{C}$ (Manning et al., 2008). Après sélection des termes, seuls les 10 000 termes les plus porteurs d'information sur les classes ont alors été retenus.

2.4 Similarité cosinus

La quatrième approche, seulement utilisée pour la prédiction du type de plat, repose elle aussi sur une représentation vectorielle des documents : à chaque recette r ,

nous faisons correspondre un vecteur v_r dont les composantes $w_r(t)$ dans l'espace du lexique sont données, pour chaque terme t présent dans le corps de la recette, par :

$$\begin{aligned} w_r(t) &= tf(t).idf(t).G(t) \\ &= tf(t).idf(t). \sum_{c \in \mathbb{C}} \mathbb{P}^2(c|t) \\ &= tf(t).idf(t). \sum_{c \in \mathbb{C}} \left(\frac{df_c(t)}{df_{\mathbb{T}}(t)} \right)^2 \end{aligned} \quad (4)$$

où $G(t)$ désigne l'indice de *Gini*, $df_{\mathbb{T}}(t)$ est le nombre de recettes du corpus d'apprentissage $\mathbb{T} \subset \mathbb{X}$ contenant le terme t et $df_c(t)$ correspond au nombre de recettes du corpus d'apprentissage annotées selon le type $c \in \mathbb{C}$.

A chaque classe $c \in \mathbb{C}$, nous faisons par ailleurs correspondre un vecteur dont les composantes $w_c(t)$ sont données dans l'espace du lexique par :

$$w_c(t) = df_c(t).idf(t).G(t) \quad (5)$$

où $df_c(t)$ désigne le nombre de recettes annotées selon le type de plat c où le terme t apparaît.

Une mesure de similarité $s(r, c)$ entre une recette $r \in \mathbb{X}$ et un type de plat $c \in \mathbb{C}$ est alors donnée par le cosinus de l'angle formé par les vecteurs v_r et v_c :

$$s(r, c) = \cos(\widehat{v_r, v_c}) = \frac{\sum_{t \in r \cap c} w_r(t).w_c(t)}{\sqrt{\sum_t w_r(t)^2.w_c(t)^2}} \quad (6)$$

Dans ces conditions la classe $\gamma(c)$ attribuée à une recette r est celle qui maximise la mesure de similarité cosinus :

$$\gamma(r) = \arg \max_{c \in \mathbb{C}} (s(r, c)) \quad (7)$$

Le vocabulaire utilisé lors de la phase de vectorisation était constitué des seuls termes t dont l'indice de Gini $G(t)$ était au moins égal à 0.45.

Ce seuil a été fixé empiriquement par maximisation du macro F-score sur un corpus de développement constitué de 3 864 des 13 864 recettes du corpus d'apprentissage.

2.5 Fusion des systèmes

Pour chacune des tâches de classification, plusieurs méthodes ont donc été appliquées : trois pour la tâche 1, quatre pour la tâche 2.

Pour chaque couple $(r, c) \in \mathbb{X} \times \mathbb{C}$ associant une recette du corpus à une classe, nous disposons donc d'un score $s_i(r, c)$ pour chacune des méthodes de classification (avec $i = 1, \dots, 3$ ou $i = 1, \dots, 4$ selon la tâche).

2.5.1 Normalisation des résultats par méthode

Les scores produits en sortie des différents systèmes de classification ont d'abord été normalisés de manière à ce que la somme des prédictions sur l'ensemble des classes pour une méthode donnée soit égale à 1. En notant $n_i(r, c)$ le score normalisé obtenu selon la i -ème méthode pour le couple $(r, c) \in \mathbb{X} \times \mathbb{C}$, nous avons donc :

$$n_i(r, c) = \frac{s_i(r, c)}{\sum_{j=1}^{|\mathbb{C}|} s_i(r, c_j)} \quad (8)$$

où c_j désigne la j -ème classe ($j = 1, \dots, 4$ ou $j = 1, \dots, 3$ selon la tâche de classification).

Deux méthodes de fusion ont été appliquées après normalisation des scores $s_i(r, c)$.

2.5.2 Combinaison linéaire des scores

La première méthode de fusion a consisté, pour une recette donnée r , à départager les classes candidates par combinaison linéaire des scores obtenus selon les différentes méthodes de classification. En notant $\gamma(r)$ la classe conjecturée, on a donc :

$$\gamma(r) = \arg \max_{c \in \mathbb{C}} \left(\sum_{i=1}^m n_i(r, c) \right) \quad (9)$$

où m correspond au nombre de méthodes appliquées pour la tâche considérée.

2.5.3 Méthode ELECTRE

Issue du domaine de l'optimisation multi-objectif discrète, la méthode ELECTRE (Roy, 1991) consiste à définir sur l'ensemble \mathbb{C} des classes candidates une relation $\mathcal{S} \subset \mathbb{C} \times \mathbb{C}$ dite de surclassement. De manière informelle, nous pouvons dire qu'une classe c surclasse une classe c' si elle la domine sur un nombre « important » de méthodes et si, sur les éventuelles méthodes restantes où elle est dominée par c' , elle ne l'est pas au-delà d'un certain seuil fixé pour chaque méthode. Plus formellement :

$$c \mathcal{S} c' \iff \begin{cases} \text{conc}(c, c') \geq sc \\ v(c, c') = 0 \end{cases} \quad (10)$$

où $\text{conc}(c, c') \in [0, 1]$ désigne l'indice de concordance entre les classes candidates c et c' , $sc \in [0, 1]$ un seuil dit de concordance fixé empiriquement et $v(c, c') \in \{0, 1\}$ un indice binaire dit de veto.

L'indice de concordance $\text{conc}(c, c')$ évalue le taux de méthodes (éventuellement pondérées) selon lesquelles c domine c' :

$$\text{conc}(c, c') = \frac{\sum_{i \in M(c, c')} P_i}{\sum_{i \in M} P_i} \quad (11)$$

où M désigne l'ensemble des méthodes, $M(c, c')$ l'ensemble des méthodes sur lesquelles c domine (non strictement) c' et p_i le poids accordé à chaque méthode $i \in M$.

L'indice binaire $v(c, c')$ (à valeurs dans l'ensemble $\{0, 1\}$) fixe la valeur du veto de c' vers c selon les modalités suivantes (en notant $n_i(r, c)$ le score obtenu selon la i -ème méthode par le couple (r, c) associant une recette à une classe) :

$$v(c, c') = \begin{cases} 1 & \iff \exists i \in M : n_i(r, c') > n_i(r, c) \text{ et } n_i(r, c') - n_i(r, c) \geq v_i \\ 0 & \text{sinon} \end{cases} \quad (12)$$

où $v(i) \in [0, 1]$ est la valeur de veto associée à la i -ème méthode.

Les valeurs des différents paramètres ont été fixées empiriquement à partir du corpus de développement par maximisation des métriques appliquées pour évaluer la classification.

Les poids p_i associés aux différentes méthodes ont tous été fixés à $p_i = 1$. Le seuil de concordance a été fixé à $sc = 0,7$ pour la tâche 1 et $sc = 0,6$ pour la tâche 2 ; enfin, la valeur de veto v_i a été fixée à $v_i = 0,5$ pour chacune des méthodes et quelle que soit la tâche considérée.

Le noyau de la relation $\mathcal{S} \subset \mathbb{C} \times \mathbb{C}$ est alors constitué des éventuelles classes candidates non surclassées.

Dans le cas d'un noyau *singleton*, la classe conjecturée pour la recette est l'unique classe élément du noyau.

Dans le cas de noyau vide ou de noyau formé de plusieurs classes concurrentes, c'est la meilleure classe candidate selon la méthode de combinaison linéaire qui a été retenue.

2.6 Extraction d'ingrédients

La tâche 4 est une tâche pilote d'extraction d'information. La difficulté est multiple, car les auteurs des recettes n'utilisent pas tous les ingrédients, ou ils ont introduit librement des ingrédients non listés. Par exemple, il n'est pas rare d'avoir des passages comme « mélanger énergiquement tous les ingrédients » (recette 27174), qui ne permettent pas d'extraire directement les ingrédients.

Nous avons attaqué ce problème en utilisant deux approches : l'une à base de règles et l'autre probabiliste. L'idée de base pour les règles a été la suivante : soit A l'ensemble de mots de la recette i ; soit B l'ensemble de mots de la liste DEFT (références d'apprentissage ou *gold-standard* d'évaluation). Éliminer tous les mots de A qui ne sont pas dans B . Cela produit une liste d'ingrédients candidats L_c .

Cette liste peut contenir des termes génériques comme « VIANDE », « FROMAGE » ou « POISSON », présents dans le texte de la recette. Afin de mieux identifier ce terme, nous avons employé une approche probabiliste naïve. Nous avons calculé la probabilité d'avoir un certain type de viande x , étant donné une liste d'ingrédients $L_c = (L_1, L_2, \dots, L_n)$:

$$p(x|L_c) = P(x \cap L_c) / p(L_c); \quad x = \{\text{VIANDE, FROMAGE, POISSON}\} \quad (13)$$

Les termes ainsi identifiés, sont alors injectés dans la liste extraite L_c , ce qui produit la liste définitive L .

La liste d'ingrédients L ainsi obtenue a été utilisée dans les systèmes de *boosting* (c.f. section 2.2) dans les tâches T1 et T2.

3 Résultats et discussion

Les systèmes ont été évalués en utilisant les mesures proposées par DEFT : le F-score pour les tâches 1 et 2, et la mesure *Mean Average Precision* (MAP) de TREC pour la tâche d'extraction.

3.1 Tâches de classification

Pour la tâche de classification en niveau de difficulté (T1), nos trois runs étaient donnés par :

- Run 1 : SVMs seuls ;
- Run 2 : fusion des trois méthodes par ELECTRE ;
- Run 3 : combinaison linéaire des trois méthodes.

Pour la tâche de classification en type de plat (T2) :

- Run 1 : Modèle discriminant seul ;
- Run 2 : fusion des quatre méthodes par ELECTRE ;
- Run 3 : combinaison linéaire des quatre méthodes.

Les résultats obtenus sont récapitulés dans les tableaux 1 (tâche 1) et 2 (tâche 2).

Pour la tâche 1 et pour chacun des trois runs, nous fournissons, outre les F-scores, la distance moyenne (micro-écart) de l'hypothèse à la référence sur l'échelle des quatre niveaux de difficulté, deux niveaux contigus étant distants de 1.

T1	run 1	run 2	run 3
Micro F-score	0.5916	0.5812	0.5877
Macro F-score	0.4531	0.4531	0.4569
Distance moyenne	0.4409	0.4586	0.4465

TABLE 1 – Résultats obtenus sur le corpus de test - Tâche 1

T2	run 1	run 2	run 3
Micro F-score	0.8760	0.8860	0.8886
Macro F-score	0.8706	0.8795	0.8824

TABLE 2 – Résultats obtenus sur le corpus de test - Tâche 2

Pour la tâche 2, la combinaison des méthodes permet d'obtenir des scores systématiquement supérieurs à ceux que nous obtenons en appliquant une méthode isolée.

Pour la tâche d'évaluation du niveau de difficulté, les bénéfices retirés de l'utilisation d'une méthode de fusion sont moins nets et apparaissent dépendants de la métrique appliquée lors de l'évaluation.

Ceci a permis au LIA de se positionner dans la tâche 1 : 3ème/6 et dans la tâche 2 : 1er/5.

En ce qui concerne la tâche T2, nous croyons que l'évaluation aurait dû prendre en compte la dimension multi-labels de plusieurs recettes. Par exemple la recette 18 052 (présente dans le test) se termine par « servir tiède (aussi bon chaud que froid ou réchauffé), en entrée ou en plat » est annotée uniquement comme une entrée. Si un système produit l'étiquette *plat principal*, faut-il pour autant compter cela comme une erreur ?

Relevons aussi l'aspect subjectif de la tâche T1. En effet, dire qu'une recette est facile ou difficile sans tenir compte qui en est l'auteur c'est faire fi de son niveau d'expertise!!!

Nous émettons également l'idée que les macro-mesures devraient être privilégiées aux micro-mesures car ces dernières pousseraient à mettre en œuvre des stratégies négligeant les classes minoritaires. Nous pensons donc que cette façon d'évaluer serait plus adaptée à une utilisation orientée application.

3.2 Tâche d'extraction

L'évaluation de cette tâche est assez subjective, malgré son caractère d'extraction d'information. Par exemple, dans l'ensemble d'apprentissage, la recette 10 514 dit : « 40 cl d'Apremont ou autre blanc de Savoie sec (facultatif, mais donne plus de goût) ». Les références DEFT indiquent **mais** comme ingrédient de la tartiflette. Evidemment nous n'avons pas incorporé l'information des ingrédients d'apprentissage, à la différence de la méthode DEFT ayant généré le *gold-standard*.

Nous avons obtenu un score de MAP=0.6364 en mesure *qrel* dans le corpus d'apprentissage, et dans le classement officiel un score de MAP=0.6287 dans le corpus de test. Cela a positionné l'équipe du LIA au rang 3ème/5. Bien que standard, la mesure MAP pourrait avoir intégré les accents dans l'évaluation des ingrédients extraits : cela aurait permis de lever des ambiguïtés qui ont été inutilement introduites par la désaccentuation.

Nous pensons que les références fournies par DEFT dans l'ensemble d'apprentissage, ont contribué à rendre la tâche floue. En effet, cela n'a pas de sens d'extraire *papier sulfurisé* ou *couteau* (extraite à partir de *pointe de couteau*) comme ingrédients. De même, l'ingrédient *maïs*, incorrectement extrait des recettes d'apprentissage comportant l'expression *mais...*, n'est pas apparu dans les recettes de test comportant ladite expression.

3.3 Comparaison versus les humains

Nous avons mis en place une petite expérience impliquant des êtres humains. Pour la tâche T4, nous avons demandé à 7 annotateurs (experts ou non en matière culinaire), d'effectuer la tâche DEFT. Ceci nous a permis de positionner nos systèmes

par rapport aux performances des personnes. Puisque les personnes ont des connaissances extra-linguistiques, les juges humains ont eu comme seule consigne celle de ne pas consulter des ressources externes (par exemple des ressources électroniques ou des livres de cuisine) pour classer les recettes ou pour extraire les ingrédients. Le tableau 3 montre la moyenne MAP pour chaque annotateur A_i sur 50 recettes de la tâche T4, choisies au hasard. La moyenne générale pour les personnes est de $MAP=0.5433$ et pour le système S de $MAP=0.6570$. Ceci montre que dans ce sous-ensemble, nos systèmes sont au-dessus des performances atteintes par les humains.

A1	A2	A3	A4	A5	A6	A7	S
0.5333	0.6029	0.5271	0.5880	0.5313	0.4679	0.5529	0.6570

TABLE 3 – Performance des personnes sur la tâche T4.

Il est clair que, même les humains ayant une expertise culinaire, ont obtenu de piètres notes dans cette tâche. Cela s'explique par la mauvaise qualité du *gold standard* fourni par DEFT. Par exemple, la crème ou le café (cuillère à café, recette 90 806) semblent subir une extraction de nature aléatoire... D'ailleurs, nous avons souvent été surpris par les ingrédients de référence de DEFT utilisés dans telle ou telle recette, qui ne peuvent en aucun cas avoir été employés de façon conjointe. Par exemple, dans la recette 64 761 (dessert), la référence DEFT liste parmi les ingrédients « moules » et « caramel » ! Évidemment il s'agit de l'ustensile *moule à charlotte*. Malheureusement il s'agit de problèmes récurrents. Sans parler de l'eau ou de la pâte (déclinée comme pâté, pâtes fraîches, brisée, sablée, etc.) qui semblent être détectées aléatoirement. Pour répondre au défi, nous avons développé plusieurs systèmes d'extraction d'ingrédients qui ont dû être modifiés (souvent à l'encontre de ce qui nous semblait logique) afin de rendre une sortie proche de celle des références. Ceci revient en fin de compte à modéliser la machine d'extraction de DEFT. Les références étant de qualité discutable, est-il intéressant de chercher à reproduire la sortie d'une telle machine ?

4 Conclusions

Malgré la quantité d'erreurs présentes dans le corpus d'apprentissage, nos algorithmes ont conduit à des résultats intéressants. Nous voulons ajouter quelques commentaires par rapport à ce défi. Nous avons observé que l'auto-évaluation du niveau de difficulté des recettes par ceux qui les publient est un exercice très subjectif, assez souvent sujet à caution. En ce qui concerne le type de plat, l'évolution du mode de vie rend assez floue la distinction entre plat principal et hors d'œuvre. Les tartes salées, les quiches ou les tartines sont de plus en plus considérées comme des plats « de résistance ». Ceci conduit souvent les internautes à exprimer de façon explicite leur indécision quant à la catégorie à retenir. Ces observations n'ont pas suffi à atténuer notre perplexité lorsque nous avons constaté que certaines méthodes à base de « cuisine » algorithmique réussissaient à faire mieux que des approches sophistiquées. Pour autant, nous n'avons pas cédé à la tentation de les retenir pour en faire une soumission !

Remerciements

Nous remercions les annotateurs qui ont bien voulu nous aider dans cette tâche ! Patricia Velázquez, Sulan Wong, Mariana Tello, Alejandro Molina et Grégoire Moreau. Nous tenons aussi à remercier les organisateurs de la campagne d'évaluation, sans qui nous n'aurions pu participer à ce défi.

Références

- Favre, B., Hakkani-Tür, D., and Cuendet, S. (2007). Icsiboost. <http://code.google.com/p/icsiboost>.
- Manning, C. D., Raghavan, P., and Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA.
- Roy, B. (1991). The outranking approach and the foundations of ELECTRE methods. *Theory and Decision*, 31 :49–73.
- Schapire, R. E. (2003). The Boosting Approach to Machine Learning : An Overview. In *Nonlinear Estimation and Classification*.
- Schapire, R. E. and Singer, Y. (2000). BoosTexter : A boosting-based system for text categorization. In *Machine Learning*, volume 39, pages 135–168.
- Torres-Moreno, J., El-Bèze, M., Bellot, P., and F., B. (2011). *Peut-on voir la détection d'opinions comme un problème de classification thématique ?*, chapter 9. Hermes Lavoisier.